

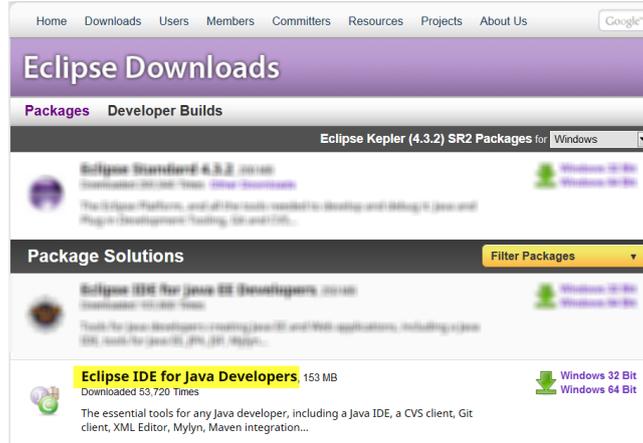
TOBB ETU HADOOP - IBM BigInsights Örnek Uygulama

İrfan Bahadır KATIPOĞLU*

6 Mart 2014

1 Çalışma Ortamının Edinilmesi

Eclipse çalışma ortamının “*Eclipse IDE for Java EE Developers*” sürümünü indiriniz.¹



Şekil 1: Eclipse Download Page

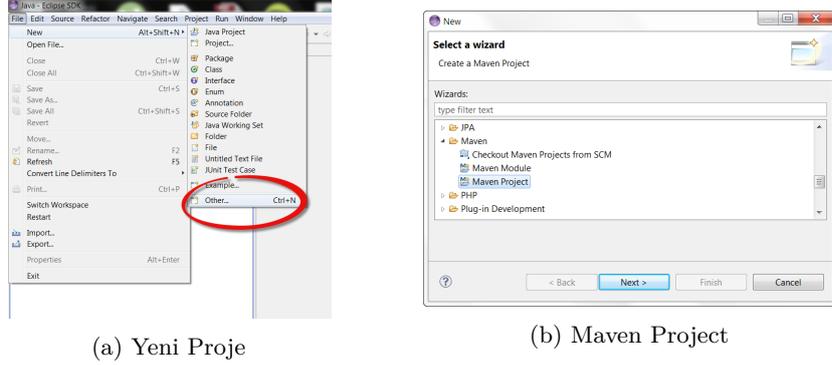
*ibkatipoglu@etu.edu.tr — bahadir@bahadir.me

¹<http://www.eclipse.org/downloads/>

2 Projenin Geliştirilmesi

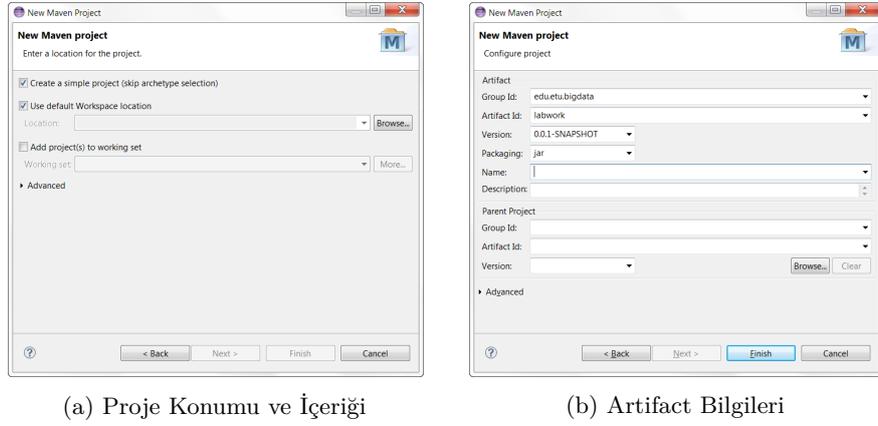
2.1 Maven Projesinin Hazırlanması

File → New → Other (Ctrl+N) ile gelen penceredeki Maven bölümünden *Maven Project*'i seçiniz.



Şekil 2: Maven Projesi

Sihirbazın sonraki ekranlarını şekildeki gibi tamamlayınız.



Şekil 3: Maven Proje Sihirbazı

Sihirbaz tamamlandıktan sonra proje iskeleti oluşacaktır. Oluşan projedeki pom.xml dosyasını açarak Kod-1ⁱ ile aynı olacak şekilde düzenleyiniz.

```
1 <project xmlns="http://maven.apache.org/POM/4.0.0" ↵  
2   ↵ xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"  
3   xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 ↵  
4     ↵ http://maven.apache.org/xsd/maven-4.0.0.xsd">  
5 <modelVersion>4.0.0</modelVersion>  
6 <groupId>edu.etu.bigdata</groupId>  
7 <artifactId>labwork</artifactId>  
8 <version>0.0.1-SNAPSHOT</version>  
9 <properties>
```

ⁱBu kodu <https://gist.github.com/bahadrix/9366787> adresinde de bulabilirsiniz.

```

9 <finalName>hadoop-labwork</finalName>
10 <mainClass>WordCount</mainClass>
11 </properties>
12
13 <dependencies>
14   <dependency> <!-- Galazy -->
15     <groupId>org.apache.hadoop</groupId>
16     <artifactId>hadoop-client</artifactId>
17     <version>2.2.0</version>
18     </dependency>
19
20   <dependency> <!-- Unit Testing ability for map reduce jobs -->
21     <groupId>org.apache.mrunit</groupId>
22     <artifactId>mrunit</artifactId>
23     <version>0.9.0-incubating</version>
24     <classifier>hadoop1</classifier>
25     </dependency>
26
27   <dependency> <!-- Apache's Common Utilities -->
28     <groupId>org.apache.commons</groupId>
29     <artifactId>commons-io</artifactId>
30     <version>1.3.2</version>
31     </dependency>
32 </dependencies>
33
34 <build>
35 <plugins>
36
37 <plugin><!-- needed for debugging in IntelliJ IDEA -->
38 <groupId>org.apache.maven.plugins</groupId>
39 <artifactId>maven-surefire-plugin</artifactId>
40 <version>2.14</version>
41 <configuration>
42   <forkMode>never</forkMode>
43 </configuration>
44 </plugin>
45
46 <plugin><!-- using shade plugin for uber jar -->
47 <groupId>org.apache.maven.plugins</groupId>
48 <artifactId>maven-shade-plugin</artifactId>
49 <version>2.2</version>
50 <executions>
51   <execution>
52     <id>main</id>
53     <phase>package</phase>
54     <goals>
55       <goal>shade</goal>
56     </goals>
57     <configuration>
58       <finalName>${finalName}</finalName>
59       <transformers>
60         <transformer
61           implementation="org.apache.maven.plugins.shade.resource.ManifestResourceTransformer">
62           <mainClass>${mainClass}</mainClass>
63         </transformer>
64       </transformers>
65       <!-- Exclude with dependencies -->
66       <minimizeJar>true</minimizeJar>
67       <artifactSet>
68         <excludes>
69           <!-- Some libs is already included in running ↴
70             ↵ machine's classpath,
71             we exclude them -->
72           <exclude>org.apache.hadoop</exclude>
73           <exclude>junit</exclude>
74           <exclude>org.apache.mrunit</exclude>
75           <exclude>log4j</exclude>
76           <exclude>org.xerial.snappy</exclude>
77           <exclude>org.mockito</exclude>
78           <exclude>org.codehaus.jackson</exclude>
79           <exclude>org.objenesis</exclude>
80           <exclude>org.apache.avro</exclude>
81           <exclude>com.google.protobuf</exclude>
82           <exclude>com.google.common</exclude>

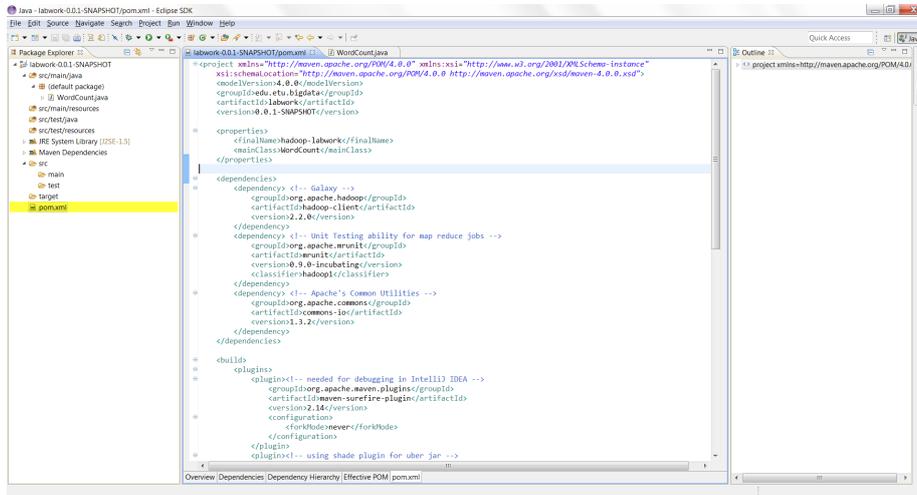
```

```

82         <exclude>org.slf4j</exclude>
83     </excludes>
84 </artifactSet>
85
86 </configuration>
87 </execution>
88 </executions>
89 </plugin>
90
91 </plugins>
92 </build>
93
94 </project>

```

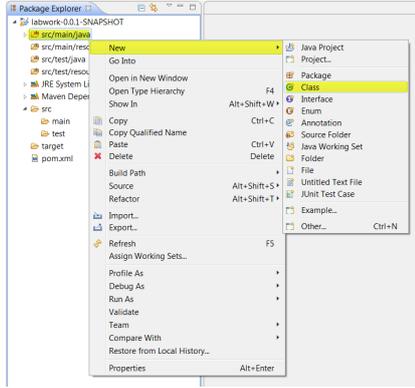
Kod 1: POM Dosyası



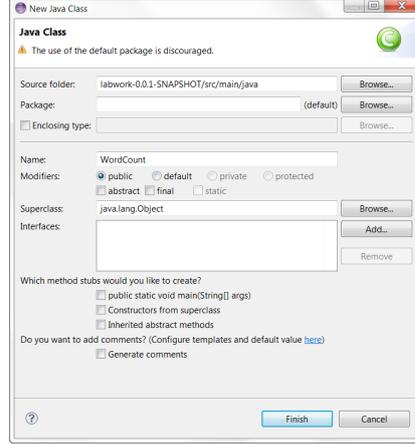
Şekil 4: Maven Proje İskeleti

2.1.1 Sınıfların Kodlanması

`src/main/java` dizinine sağ tıklayıp New→Class seçeneği ile şekildeki gibi bir `WordCount` sınıfı oluşturunuz.



(a) Yeni Sınıf Oluşturma



(b) Sınıf Özellikleri

Şekil 5: WordCount Sınıfının Oluşturulması

Bu sınıfı açarak Kod-2¹ ile aynı olacak şekilde düzenleyiniz.

```
1  /**
2   * ETU Hadoop - BigInsights
3   * Sample WordCount Class
4   * @author Bahadır Kaitpoglu
5   */
6
7  import org.apache.hadoop.conf.Configuration;
8  import org.apache.hadoop.conf.Configured;
9  import org.apache.hadoop.fs.Path;
10 import org.apache.hadoop.io.IntWritable;
11 import org.apache.hadoop.io.LongWritable;
12 import org.apache.hadoop.io.Text;
13 import org.apache.hadoop.mapred.*;
14 import org.apache.hadoop.util.Tool;
15 import org.apache.hadoop.util.ToolRunner;
16
17 import java.io.IOException;
18 import java.util.Iterator;
19 import java.util.StringTokenizer;
20
21 public class WordCount extends Configured implements Tool {
22
23     /**
24     * Mapper class
25     */
26     /**
27     static class Mappa extends MapReduceBase implements
28     Mapper<LongWritable, Text, Text, IntWritable> {
29
30         private Text word = new Text();
31         private final static IntWritable ONE = new IntWritable(1);
32
33         public void map(LongWritable key, Text value,
34             OutputCollector<Text, IntWritable> collector, Reporter arg3)
35             throws IOException {
36
```

¹Kodu <https://gist.github.com/bahadrix/9370719> adresinden de indirebilirsiniz.

```

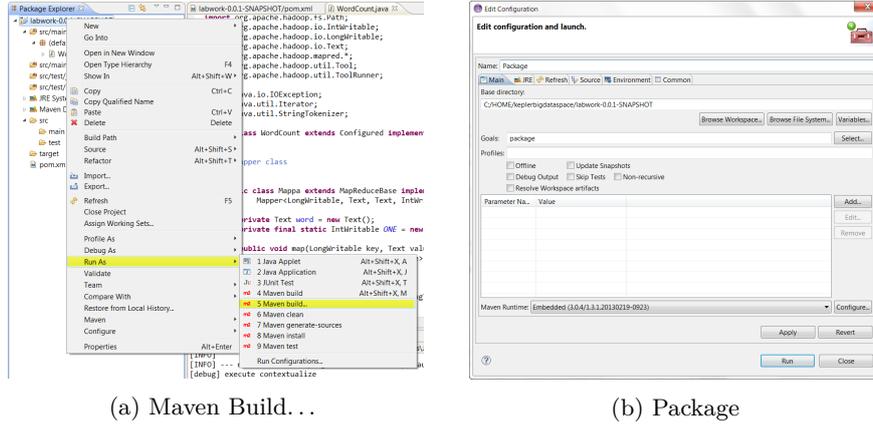
37     String line = value.toString();
38     StringTokenizer tokenizer = new StringTokenizer(line);
39     while (tokenizer.hasMoreTokens()) {
40         word.set(tokenizer.nextToken());
41         collector.collect(word, ONE);
42     }
43
44     }
45
46 }
47
48 /**
49  * Reducer Class
50  *
51  */
52 static class Reducca extends MapReduceBase implements
53     Reducer<Text, IntWritable, Text, IntWritable> {
54
55     public void reduce(Text key, Iterator<IntWritable> values,
56         OutputCollector<Text, IntWritable> collector, Reporter arg3)
57         throws IOException {
58         int sum = 0;
59         while (values.hasNext()) {
60             sum += values.next().get();
61         }
62         collector.collect(key, new IntWritable(sum));
63     }
64 }
65
66 }
67
68 /**
69  * Entry point
70  *
71  * @param args
72  * @throws Exception
73  */
74 public static void main(String[] args) throws Exception {
75     System.exit(ToolRunner.run(new Configuration(), new WordCount(),
76         args));
77 }
78
79 /**
80  * Run the job
81  */
82 public int run(String[] args) throws Exception {
83     Configuration conf = getConf();
84     JobConf job = new JobConf(conf, WordCount.class);
85
86     Path pIn = new Path(args[0]);
87     Path pOut = new Path(args[1]);
88     FileInputFormat.setInputPaths(job, pIn);
89     FileOutputFormat.setOutputPath(job, pOut);
90
91     job.setJobName("WordCount");
92     job.setMapperClass(Mappa.class);
93     job.setReducerClass(Reducca.class);
94
95     job.setOutputKeyClass(Text.class);
96     job.setOutputValueClass(IntWritable.class);
97
98     JobClient.runJob(job);
99
100     return 0;
101 }
102
103 }

```

Kod 2: WordCount Simfi

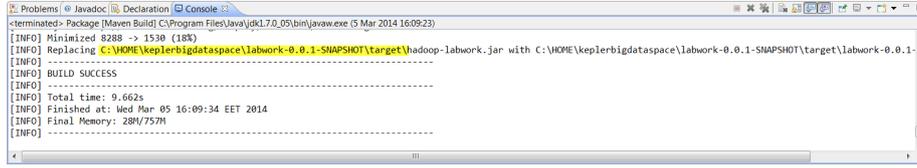
3 Derleme

Projenin ana dizinine sağ tıklayıp *Run as*→*Maven Build...* seçeneğini seçiniz. Gelen ekrandaki *goals* ve *package* alanlarını şekildeki gibi doldurup Run butonuna tıklayınız¹.



Şekil 6: Maven Package Konfigurasyonunun Oluşturulması

Maven paketlemeyi başardığında aşağıdaki gibi bir konsol çıktısı vermeli. Şekilde işaretlenmiş olan adres Hadoop ile kullanılmak üzere oluşturulmuş olan JAR dosyasının bilgisayarındaki konumunu gösteriyor. Bu konumu kopyalıyoruz.



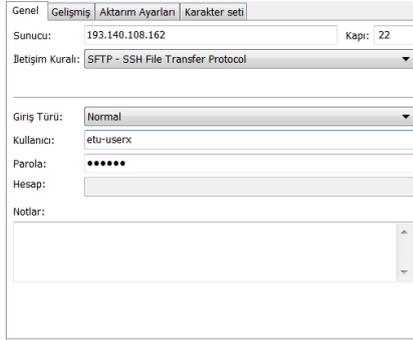
Şekil 7: Maven Paketleme İşlemi Konsol Çıktısı

4 Yükleme

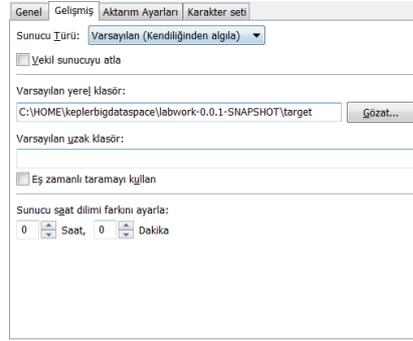
4.1 SFTP İşlemleri

FileZilla adlı programı https://filezilla-project.org/download.php?show_all=1 adresinden indirerek kurunuz. Programı çalıştırıp Ctr + S tuşlarına basarak gelen SiteManager ekranından New Site butonuna tıklayarak yapılandırmasını aşağıdaki şekilde gösterildiği gibi yapınız. Yerel dizin olarak jar dosyasının bulunduğu dizinin adresini yazınız.

¹Daha sonradan projeyi tekrar paketlemek için Run→Run History→Package menüsünü kullanabilirsiniz.



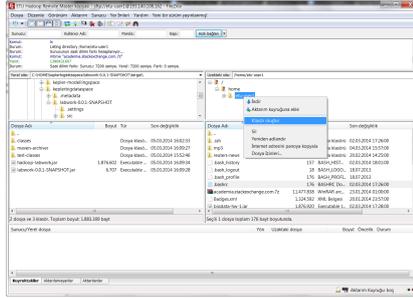
(a) Sunucu ayarları



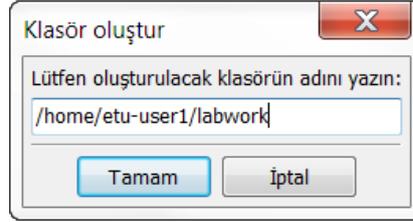
(b) Yerel dizin ayarı

Şekil 8: FileZilla SFTP erişiminin oluşturulması

Bağlan butonuna basarak MasterNode'a bağlanınız. Sol taraf bağlanan bilgisayarın yerel dizinini, sağ taraf sunucu tarafını göstermektedir. Bağlantı kurulduktan sonra aşağıda gösterildiği gibi sunucudaki dizininize *labwork* adlı bir alt dizin ekleyiniz.



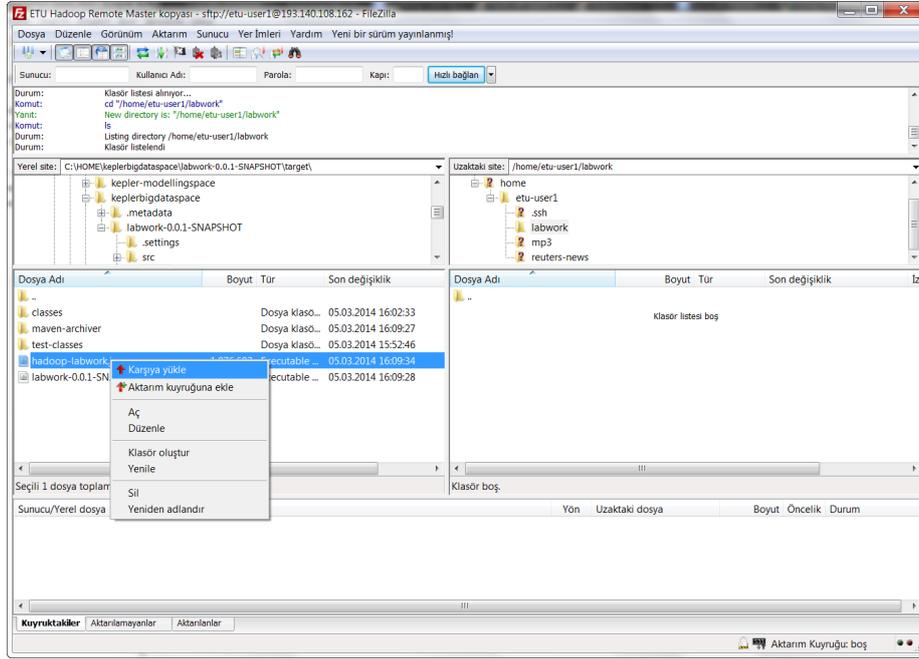
(a) Kullanıcı dizini altında yeni klasör



(b) labwork

Şekil 9: Çalışma dizininin oluşturulması

Daha sonra oluşturmuş olduğunuz jar dosyasını bu yeni dizine kopyalayınız.



Şekil 10: JAR Dosyasının Transferi

Bu işlemden sonra ödevde verilmiş olan reuters-news.rar dosyasını indiriniz.

TOBB University of Economics and Technology - Spring 2014

BİL 401/501: Big Data

Syllabus

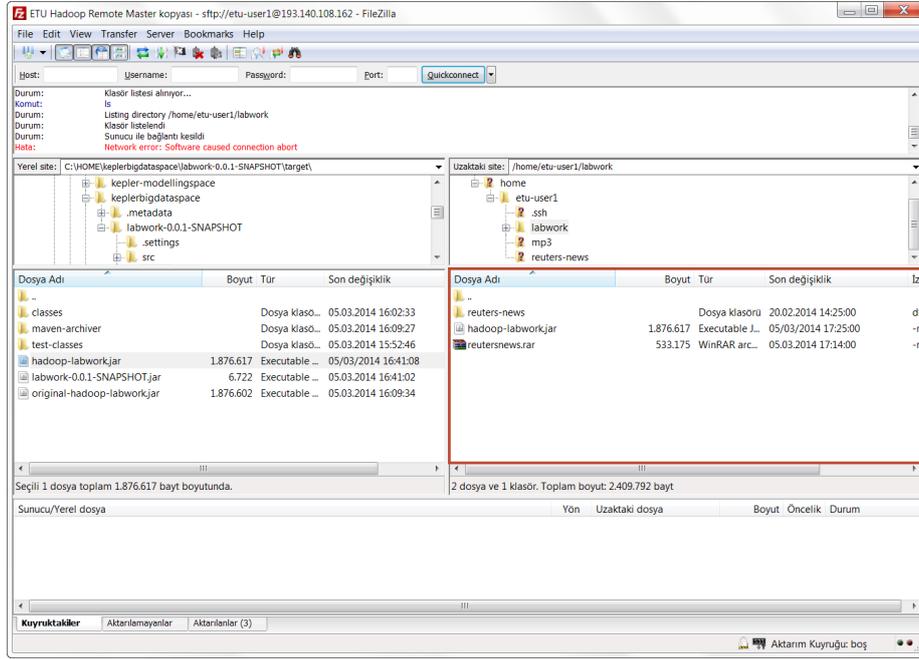
Course Information Staff Resources

Homework

Homework	Due Date
reutersnews.rar	Mar 6, 2014
hw1.rar	Mar 6, 2014

Şekil 11: Piazza - ReutersNews

İndirdiğiniz dosyayı Şekil 12 ile gösterilen alana sürükleyip bırakarak sunucudaki aynı dizine transferini sağlayınız.



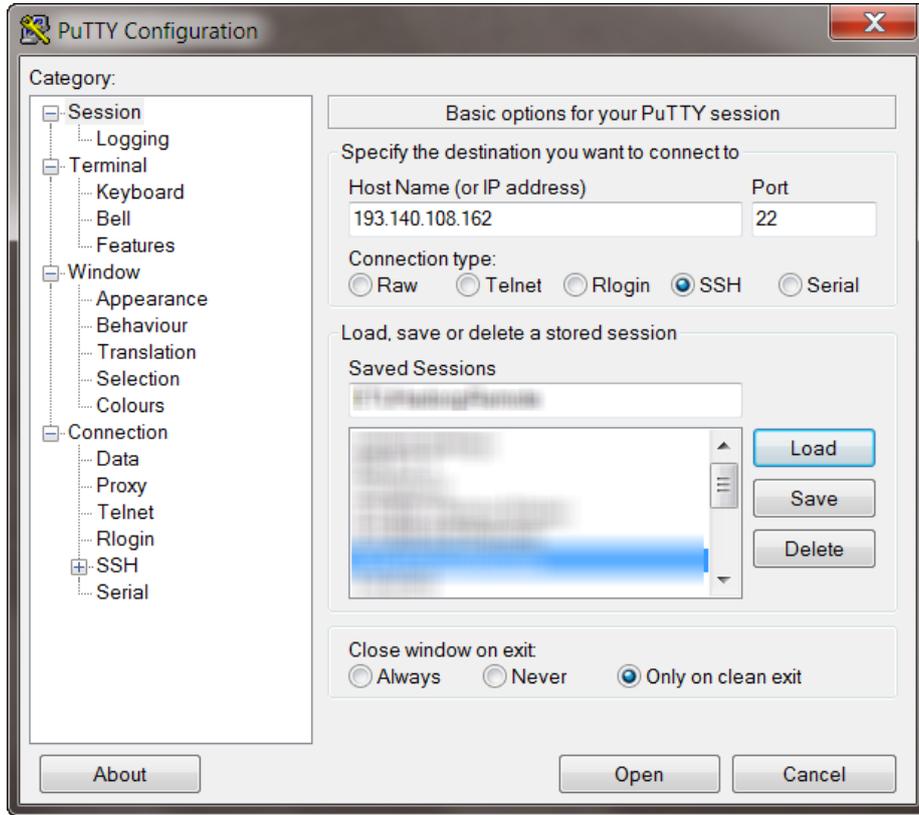
Şekil 12: Veri Setinin Transferi

5 Uygulamanın Hadoop İle Yürütülmesi

5.1 SSH Bağlantısının Sağlanması

Windows kullanıcıları Putty ile şekilde gösterildiği gibi sunucuya bağlanabilirler. Linux tabanlı sistemlerde aşağıdaki gibi sunuya SSH bağlantısı kurulabilir.

```
ssh 193.140.108.162
```



Şekil 13: Putty yapılandırma ekranı

5.2 SSH İşlemleri

```
1 $ cd labwork
2
3 # Check our uploaded files
4 $ ls
5 hadoop-labwork.jar reutersnews.rar
6
7 # Extract dataset
8 $ unrar x reutersnews.rar
9
10 # Create folders for dataset
11 $ hadoop dfs -mkdir labwork/input
12
13 # Put files to the HDFS.
14 # time is optional for measuring time
15 $ time hadoop dfs -put reuters-news/*.txt labwork/input
16 real    0m6.506s
17 user    0m5.870s
18 sys    0m0.702s
19
20 # Run!
21 $ hadoop jar hadoop-labwork.jar labwork/input labwork/output
22
23 # Check the output
24 $ hadoop dfs -ls labwork/output
25 Found 3 items
26 -rw-r--r--  3 ... /user/etu-user1/labwork/output/_SUCCESS
27 drwxr-xr-x  - ... /user/etu-user1/labwork/output/_logs
28 -rw-r--r--  3 ... /user/etu-user1/labwork/output/part-00000
29
30 # Print out the result
31 $ hadoop dfs -cat labwork/output/*0
32 (...)
33 zero      4
34 zero,    1
35 zero-coupon      3
36 zestril,        2
37 zinc           1
38 zone           4
39 zone,          1
40 zones,         1
```

Kod 3: Hadoop Görevinin İşletilmesi

Çalışmakta olan ve geçmiş Job'lara ait durumu BigInsights Web Console'dan Application Status bölümünden izleyebilirsiniz.